

Entropy and long-range correlations in random symbolic sequences

S. S. Melnik and O. V. Usatenko

*A. Ya. Usikov Institute for Radiophysics and Electronics
Ukrainian Academy of Science, 12 Proskura Street, 61805 Kharkov, Ukraine*

The goal of this paper is to develop an estimate for the entropy of random long-range correlated symbolic sequences with elements belonging to a finite alphabet. As a plausible model, we use the high-order additive stationary ergodic Markov chain. Supposing that the correlations between random elements of the chain are weak we express the differential entropy of the sequence by means of the symbolic pair correlation function. We also examine an algorithm for estimating the differential entropy of *finite* symbolic sequences. We show that the entropy contains two contributions, the correlation and fluctuation ones. The obtained analytical results are used for numerical evaluation of the entropy of written English texts and DNA nucleotide sequences. The developed theory opens the way for constructing a more consistent and sophisticated approach to describe the systems with strong short- and weak long-range correlations.

PACS numbers: 05.40.-a, 87.10+e

I. INTRODUCTION

Our world is complex, chaotic and correlated. The most peculiar manifestations of this concept are human and animal communication, written texts of natural languages, DNA and protein sequences, data flows in computer networks, stock indexes, solar activity, weather, etc. For this reason, systems with long-range interactions (and/or sequences with long-range memory) and natural sequences with non-trivial information content have been the focus of a large number of studies in different fields of science for the past several decades. The unflagging interest in the systems with correlated fluctuations is also explained by the specific properties they demonstrate and their prospective applications as a creative tool for designing the devices and appliances with random components in their structure (different wave-filters, diffraction gratings, artificial materials, antennas, converters, delay lines, etc. [1]).

Random sequences with *finite number of states* exist as natural sequences (DNA or natural language texts) or arise as a result of coarse-grained mapping of the evolution of the chaotic dynamical system into a string of symbols [2, 3]. Such random sequences are the subject of study of the algorithmic (Kolmogorov-Solomonoff-Chaitin) complexity, artificial intellect, information theory, compressibility of digital data, statistical inference problem, computability and have many application aspects mentioned above.

There are many methods for describing complex dynamical systems and random sequences connected with them: fractal dimensions, multi-point probability distribution functions, correlation functions, and many others. One of the most convenient characteristics serving to the purpose of studying complex dynamics is entropy [4, 5]. Being a measure of the information content and redundancy in a sequence of data, it is a powerful and popular tool in examination of complexity phenomena. Among fields of science where the notion of entropy is of major significance data compression [6], natural language

processing [7] and artificial intelligence [8] are the most important. The basic idea of compression is to exploit redundancy in data, expressed in terms of correlations, and transform this redundancy in compression algorithm. Recent advances in different fields of science have hinted at a deep connection between intelligence and entropy.

A standard method of understanding and describing statistical properties of a given random sequence of data requires the estimation of the joint probability function of words occurring for sufficiently large length L of words. For limited size sequences, reliable estimations can be achieved only for very small L because the number m^L (where m is the finite-alphabet length) of different words of the length L has to be much less than the total number $M - L$ of words in the whole sequence of the length M ,

$$m^L \ll M - L \simeq M. \quad (1)$$

This is the crucial point because usually the correlation lengths of natural sequences of interest is of the same order that the length of sequence. Inequality (1) cannot be fulfilled. The lengths of representative words that could estimate correctly the probability of words occurring are 4–5 for a real natural text of the length 10^6 (written on an alphabet containing 27–30 letters and symbols) or of order of 20 for a coarse-grained text represented through a binary sequence. So, long-range correlations that can exist in the sequences cannot be taken into account in such a kind of theories.

Here we present a complementary approach, which takes into account just the long-range correlations. Specifically, we sacrifice the knowledge of exact statistics of short words and take into account the weak long-range memory, which can be expressed in terms of the pair correlation function of symbols and can be found by numerical analysis of sequence nearly at the same distances as the total length of sequence.

We use the earlier developed method [9] for constructing the conditional probability function presented by means of pair correlator, which makes it possible to calculate analytically the entropy of the sequence. It should be

stressed that we suppose that the correlations are weak but not short. Which kind of memory, long- or short-range, is more important depends on the intrinsic correlation properties of the sequence under study.

The scope of the paper is as follows. First, supposing that the correlations between symbols in the sequence are weak, we represent the differential entropy in terms of the conditional probability function of the Markov chain and express the entropy as the sum of squares of the pair correlators. Then we discuss some properties of the results obtained. Next, a fluctuation contribution to the entropy due to finiteness of random chains is examined. The application of the developed theory to literary texts and DNA sequences of nucleotides is considered. In conclusion, some remarks on directions in which the research can be progressed are presented.

This work is a generalization of our previous paper [10] devoted to the binary random sequences. We insistently recommend to a reader to see it before reading this paper.

II. ENTROPY OF THE ADDITIVE SYMBOLIC MARKOV CHAINS

Consider a semi-infinite random stationary ergodic sequence

$$\mathbb{A} = a_0, a_1, a_2, \dots \quad (2)$$

of symbols (letters) a_i taken from the finite alphabet

$$A = \{\alpha^1, \alpha^2, \dots, \alpha^m\}, \quad a_i \in A, \quad i \in \mathbb{N}_+ = \{0, 1, 2, \dots\}. \quad (3)$$

We use the notation a_i to indicate a position of the symbol a in the chain and the notation α^k to stress the value of the symbol $a \in A$.

We suppose that the symbolic sequence \mathbb{A} is the *high-order Markov chain* [11–15]. Such sequences are also referred to as the multi- or the N -step [16–18] Markov's chains. The sequence \mathbb{A} is the N -step Markov's chain if it possesses the following property: the probability of symbol a_i to have a certain value $\alpha^k \in A$ under condition that *all* previous symbols are given depends only on N previous symbols,

$$\begin{aligned} P(a_i = \alpha^k | \dots, a_{i-2}, a_{i-1}) \\ = P(a_i = \alpha^k | a_{i-N}, \dots, a_{i-2}, a_{i-1}). \end{aligned} \quad (4)$$

Sometimes the number N is also referred to as the *order* or the *memory length* of the Markov chain. Note, definition (4) is valid for $i \geq N$; for $i < N$ we should use the well known conditions of compatibility for the conditional probability functions of lower order [19].

To estimate the differential entropy of stationary sequence \mathbb{A} of symbols a_i one could use the Shannon definition [4] for entropy per block of length L ,

$$H_L = - \sum_{a_1, \dots, a_L \in A} P(a_1^L) \log_2 P(a_1^L). \quad (5)$$

Here $P(a_1^L) = P(a_1, \dots, a_L)$ is the probability to find L -word a_1^L in the sequence; hereafter we use the more concise notation a_{i-N}^{i-1} for N -word a_{i-N}, \dots, a_{i-1} . The differential entropy, or the entropy per symbol, is given by

$$h_L = H_{L+1} - H_L. \quad (6)$$

This quantity specifies the degree of uncertainty of $(L+1)$ th symbol occurring and measures the average information per symbol if the correlations of $(L+1)$ th symbol with preceding L symbols are taken into account. The differential entropy h_L can be represented in terms of the conditional probability function $P(a_{L+1}|a_1^L)$,

$$h_L = \sum_{a_1, \dots, a_L \in A} P(a_1^L) h(a_{L+1}|a_1^L) = \overline{h(a_{L+1}|a_1^L)}, \quad (7)$$

where $h(a_{L+1}|a_1^L)$ is the amount of information contained in the $(L+1)$ th symbol of the sequence conditioned on L previous symbols,

$$h(a_{L+1}|a_1^L) = - \sum_{a_{L+1} \in A} P(a_{L+1}|a_1^L) \log_2 P(a_{L+1}|a_1^L). \quad (8)$$

The source entropy (or Shannon entropy) is the differential entropy at the asymptotic limit, $h = \lim_{L \rightarrow \infty} h_L$. This quantity measures the average information per symbol if *all* correlations, in the statistical sense, are taken into account, cf. with [20], Eq. (3).

Due to the ergodicity of stationary sequence \mathbb{A} , the average value of any function $f(a_{r_1}, a_{r_1+r_2}, \dots, a_{r_1+\dots+r_s})$ of s arguments defined on the set A of symbols is statistical (arithmetic, Cesaro's) average over the chain,

$$\begin{aligned} \overline{f}(a_{r_1}, \dots, a_{r_1+\dots+r_s}) \\ = \lim_{M \rightarrow \infty} \frac{1}{M} \sum_{i=0}^{M-1} f(a_{i+r_1}, \dots, a_{i+r_1+\dots+r_s}). \end{aligned} \quad (9)$$

Stationarity together with decay of correlations, $C_{\alpha, \beta}(r \rightarrow \infty) = 0$, see below definition (13), leads, according to the Slutsky sufficient conditions [21], to mean-ergodicity. This latter property is very useful in numerical calculations since the averaging can be done over the length of the sequence and the ensemble averaging can be avoided. Therefore, in our numerical as well as analytical calculations we always apply averaging over the length of the sequence as it is implied in Eq. (7).

If the sequence, statistical properties of which we would like to analyze, is given, the conditional probability function (CPF) of N th order can be found by a standard method (written below for subscript $i = N+1$)

$$P(a_{N+1} = \alpha^k | a_1, \dots, a_N) = \frac{P(a_1, \dots, a_N, \alpha^k)}{P(a_1, \dots, a_N)}, \quad (10)$$

where $P(a_1, \dots, a_N, \alpha^k)$ and $P(a_1, \dots, a_N)$ are the probabilities of the $(N+1)$ -subsequence $a_1, \dots, a_N, \alpha^k$ and N -subsequence a_1, \dots, a_N occurring, respectively.

The Markov chain with CPF of general form Eq. (4) is not convenient (compliant) to solve concrete problems. For this reason we introduce a simplification for the CPF. Specifically, we suppose that the symbolic Markov chain under consideration is *additive*, i.e. its conditional probability is a linear function of random variables a_k , $k = i - N, \dots, i - 1$,

$$P(a_i = \alpha | a_{i-N}^{i-1}) = p_\alpha + \sum_{r=1}^N \sum_{\beta \in A} F_{\alpha\beta}(r) [\delta(a_{i-r}, \beta) - p_\beta], \quad (11)$$

where p_α is the relative number of symbols α in the chain, or their probabilities of occurring,

$$p_\alpha = \overline{\delta(a_i, \alpha)}. \quad (12)$$

Here $\delta(.,.)$ is the Kronecker delta-symbol, playing the role of the characteristic function of the random variable a_i and converting symbols to numbers. Hereafter, we often drop the superscript k from α^k to simplify the notations.

The additivity means that the previous symbols a_{i-N}^{i-1} exert an independent effect on the probability of the symbol $a_i = \alpha$ occurring. The first term in the right-hand side of Eq. (11) is responsible for correct reproduction of statistical properties of uncorrelated sequences, the second one takes into account, and produces under generation, correlations among symbols of the random sequence. The conditional probability function in form (11) can reproduce correctly the binary (pair, two-point) correlations in the chain. Higher-order correlators and all correlation properties of higher orders are not independent anymore. We cannot control them and reproduce correctly by means of the memory function $F_{\alpha\beta}(r)$ because the latter is completely determined by the pair correlation function, see below Eq. (18).

The additive Markov chains are, in some sense, analogous to the chains described by autoregressive models [11, 22]. In Appendix A some suggestions on the form of Eq. (11) and its properties are presented.

There is a rather simple relation between the memory function $F_{\alpha\beta}(r)$ and the pair *symbolic* correlation function of the additive Markov chain. The two-point symbolic correlation function is defined as

$$C_{\alpha\beta}(r) = \overline{[\delta(a_i, \alpha) - p_\alpha][\delta(a_{i+r}, \beta) - p_\beta]}, \quad \alpha, \beta \in A. \quad (13)$$

This function possesses the following properties:

$$\begin{aligned} C_{\alpha\beta}(r) &= C_{\beta\alpha}(-r), \\ \sum_{\alpha \in A} C_{\alpha\beta}(r) &= \sum_{\beta \in A} C_{\alpha\beta}(r) = 0. \end{aligned} \quad (14)$$

Let us suppose that there exists a one-to-one correspondence $a_i \leftrightarrow \varepsilon_i$ between the letters of symbolic sequence \mathbb{A} and the numbers of numeric sequence. Then, the ordinary “numeric” correlation function

$$C_\varepsilon(r) = \overline{(\varepsilon_i - \bar{\varepsilon})(\varepsilon_{i+r} - \bar{\varepsilon})} \quad (15)$$

of the sequence of ε_i can be expressed by means of symbolic correlator

$$C_\varepsilon(r) = \sum_{\alpha, \beta \in A} \varepsilon^\alpha \varepsilon^\beta C_{\alpha\beta}(r). \quad (16)$$

Here ε^α is the numeric value of the random variable ε corresponding to the symbol α , $\sum_{\alpha \in A}$ means the summation over all possible letters of the alphabet A .

There were suggested two methods for finding $F_{\alpha\beta}(r)$ of a sequence with a known pair correlation function. The first one [9] is based on the minimization of the “distance” between the conditional probability function, containing the sought-for memory function, and the given sequence \mathbb{A} of symbols with a known correlation function,

$$Dist = \overline{[\delta(a_i, \alpha) - P(a_i = \alpha | a_{i-N}^{i-1})]^2}. \quad (17)$$

For any values of $\alpha, \beta \in A$ and $r \geq 1$ the minimization equation with respect to $F_{\alpha\beta}(r)$ yields the relationship between the correlation and memory functions,

$$C_{\alpha\beta}(r) = \sum_{r'=1}^N \sum_{\gamma \in A} C_{\alpha\gamma}(r-r') F_{\beta\gamma}(r'). \quad (18)$$

The second method for deriving Eq. (18) is a completely probabilistic straightforward calculation analogous to that used in [17].

Equation (18), despite its simplicity, can be analytically solved only in some particular cases: for one- or two-step chains, the Markov chain with a step-wise memory function and so on. To avoid the various difficulties in its solving we suppose that correlations in the sequence are weak (in amplitude, but not in length). In order to formulate this condition we introduce the *normalized* symbolic correlation function defined by

$$K_{\alpha\beta}(r) = \frac{C_{\alpha\beta}(r)}{C_{\alpha\beta}(0)}, \quad C_{\alpha\beta}(0) = p_\alpha \delta(\alpha, \beta) - p_\alpha p_\beta. \quad (19)$$

We can obtain an approximate solution for the memory function in the form of the series

$$F_{\alpha\beta}(r) = K_{\beta\alpha}(r) + \sum_{r' \neq r}^N \sum_{\gamma \in A} K_{\gamma\alpha}(r-r') K_{\beta\gamma}(r') + \dots \quad (20)$$

if we suppose the all components of the normalized correlation function with $r \neq 0$ are small with respect to $K_{\alpha\beta}(0) = 1$.

Equation (11) for the conditional probability function in the first approximation with respect to the small parameters $|K_{\alpha\beta}(r)| \ll 1$, $r \neq 0$ after neglecting the second term in Eq.(20) takes the form

$$P(a_i = \alpha | a_{i-N}^{i-1}) \simeq p_\alpha + \sum_{r=1}^N \sum_{\beta \in A} K_{\beta\alpha}(r) [\delta(a_{i-r}, \beta) - p_\beta]. \quad (21)$$

This formula provides a tool for constructing weak correlated sequences with a given pair correlation function [9]. Note that i -independence of the function $P(a_i = \alpha | a_{i-N}^{i-1})$ provides homogeneity and stationarity of the sequence under consideration; and finiteness of N together with the strict inequalities

$$0 < P(a_{i+N} = \alpha | a_{i-N}^{i-1}) < 1, \quad i \in \mathbb{N}_+ = \{0, 1, 2, \dots\} \quad (22)$$

provides, according to the Markov theorem (see, e.g., Ref. [19]), ergodicity of the sequence.

The conditional probability $P(a_i = \alpha | a_{i-L}^{i-1})$ for a word of length $L < N$ can be obtained in the first approximation in the weak correlation parameter $\Delta_\alpha(L)$ from Eqs. (11) and (21) by means of a routine probabilistic reasoning presented in Appendix B,

$$P(a_i = \alpha | a_{i-L}^{i-1}) = p_\alpha + \Delta_\alpha(L), \quad (23)$$

$$\Delta_\alpha(L) = \sum_{r=1}^L \sum_{\beta \in A} K_{\beta\alpha}(r) [\delta(a_{i-r}, \beta) - p_\beta].$$

Taking into account the weakness of correlations,

$$|\Delta_\alpha(L)| \ll 1, \quad (24)$$

we expand Eq. (8) in Taylor series up to the second order in $\Delta_\alpha(L)$, $h(a_{L+1} | a_1^L) = h_0 + (\partial h / \partial p_\alpha) \Delta_\alpha(L) + (1/2)(\partial^2 h / \partial p_\alpha^2) \Delta_\alpha^2(L)$, where the derivatives are taken at the point $P(a_i = \alpha | a_{i-L}^{i-1}) = p_\alpha$ and h_0 is the entropy of uncorrelated sequence,

$$h_0 = - \sum_{\alpha \in A} p_\alpha \log_2(p_\alpha). \quad (25)$$

Then, the differential entropy of the sequence in line with $\Delta_\alpha(L) = 0$ takes the form

$$h_L = \begin{cases} h_{L < N} = h_0 - \frac{1}{2 \ln 2} \sum_{r=1}^L \sum_{\alpha \in A} \frac{\overline{\Delta_\alpha^2(L)}}{p_\alpha}, \\ h_{L > N} = h_{L=N}. \end{cases} \quad (26)$$

If the length of block exceeds the memory length, $L > N$, the conditional probability $P(a_i = \alpha | a_{i-L}^{i-1})$ depends only on N previous symbols, see Eq. (4). Then, it is easy to show from (7) that the differential entropy remains constant at $L \geq N$. Thus, the second line in Eq. (26) is consistent with the first line because in the first approximation in the weak correlations the parameter $\Delta_\alpha(L)$ vanishes at $L > N$ together with the correlation function. The final expression, the main analytical result of the paper, for the differential entropy of a stationary ergodic weakly correlated random sequence is

$$h_L = h_0 - \frac{1}{2 \ln 2} \sum_{r=1}^L \sum_{\alpha, \beta \in A} \frac{C_{\alpha\beta}^2(r)}{p_\alpha p_\beta}. \quad (27)$$

In order to obtain this equation we used Eq. (23) and replaced the term $C_{\alpha\beta}(r' - r)$ with $C_{\alpha\beta}(0)\delta(r, r')$ when calculating the summation.

III. DISCUSSION

It follows from Eq. (27) that the additional correction to the entropy h_0 of the uncorrelated sequence is negative. This is the anticipated result – the correlations decrease the entropy. The conclusion is not sensitive to the sign of correlations: persistent correlations, $K > 0$, describing an “attraction” of the symbols of the same kind, and anti-persistent correlations, $K < 0$, corresponding to a “repulsion” between the same symbols, provide the corrections of the same negative sign. If the correlation function is constant at $1 \leq r \leq N$, the entropy is a linear decreasing function of the argument L up to the point $r = N$.

Equation (27) takes more simple form for a binary, $m = 2$, chain of symbols, which can be also considered as a numeric chain of random variables a_i with the alphabet of symbols-numbers $A = \{0; 1\}$. Let $p_1 = \bar{a}$, $p_0 = 1 - \bar{a}$. In order to calculate h_L we should calculate four symbolic correlation functions:

$$\begin{aligned} C_{11}(r) &= \overline{\delta(a_i, 1)\delta(a_{i+r}, 1)} - \bar{a}^2, \\ C_{00}(r) &= \overline{\delta(a_i, 0)\delta(a_{i+r}, 0)} - (1 - \bar{a})^2, \\ C_{01}(r) &= \overline{\delta(a_i, 0)\delta(a_{i+r}, 1)} - (1 - \bar{a})\bar{a}, \\ C_{10}(r) &= \overline{\delta(a_i, 1)\delta(a_{i+r}, 0)} - \bar{a}(1 - \bar{a}). \end{aligned} \quad (28)$$

Taking into account that $\delta(a_i, 1) = a_i$, $\delta(a_i, 0) = 1 - a_i$, we obtain

$$\begin{aligned} C_{11}(r) &= C_{00}(r) = C(r), \\ C_{01}(r) &= C_{10}(r) = -C(r). \end{aligned} \quad (29)$$

Here $C(r)$ is the ordinary numeric correlator

$$C(r) = \overline{(a_i - \bar{a})(a_{i+r} - \bar{a})}. \quad (30)$$

After simple algebra, we get

$$h_L = h_0 - \frac{1}{2 \ln 2} \sum_{r=1}^L K^2(r), \quad (31)$$

where $K(r)$ is the normalized pair correlation function of the binary sequence $K(r) = C(r)/C(0)$, the result obtained earlier in Ref. [10].

IV. FINITE RANDOM SEQUENCES

The relative numbers p_α of symbols in the chain, correlation functions and other statistical characteristics of random sequences are deterministic quantities only in the limit of their infinite lengths. It is a direct consequence of the law of large numbers. If the sequence length M is finite, the set of numbers a_1^M cannot be considered anymore as ergodic sequence. In order to restore its status we have to introduce the *ensemble* of finite sequences

$\{a_1^M\}_p, p \in \mathbb{N} = 0, 1, 2, \dots$. Yet, we would like to retain the right to examine *finite* sequences by using a single finite chain. So, for a finite chain we have to replace definition (13) of the correlation function by the following one,

$$C_{\alpha\beta, M}(r) = \frac{1}{M-r} \sum_{i=0}^{M-r-1} [\delta(a_i, \alpha) - p_\alpha] [\delta(a_{i+r}, \beta) - p_\beta],$$

$$p_\alpha = \frac{1}{M} \sum_{i=0}^{M-1} \delta(a_i, \alpha), \quad (32)$$

which coincides with Eq. (13) in the limit $M \rightarrow \infty$. Now the correlation functions and p_α are random quantities, which depend on the particular realization of the sequence a_1^M . Fluctuations of these random quantities can contribute to the entropy of finite random chains even if the correlations in the random sequence are absent. It is well known that the order of relative fluctuations of additive random quantity (as, e.g. the correlation function Eq. (32)) is $1/\sqrt{M}$.

Below we give more rigorous justification of this explanation and show its applicability to our case. Let us present the correlation function $C_M(r)$ as the sum of two components,

$$C_{\alpha\beta, M}(r) = C_{\alpha\beta}(r) + C_{\alpha\beta, f}(r), \quad r \geq 1, \quad (33)$$

where the first summand $C_{\alpha\beta}(r) = \lim_{M \rightarrow \infty} C_{\alpha\beta, M}(r)$ is the correlation function determined by Eq. (32) (in the limit $M \rightarrow \infty$) obtained by averaging over the sequence with respect to index i , enumerating the elements a_i of sequence \mathbb{A} ; and the second one, $C_{\alpha\beta, f}(r)$, is a fluctuation-dependent contribution. Function $C_{\alpha\beta}(r)$ can be also presented as the ensemble average $C_{\alpha\beta}(r) = \langle C_{\alpha\beta, M}(r) \rangle$ due to the ergodicity of the (infinite) sequence.

Now we can find a relationship between variances of $C_{\alpha\beta, M}(r)$ and $C_{\alpha\beta, f}(r)$. Taking into account Eq. (33) and the properties $\langle C_{\alpha\beta, f}(r) \rangle = 0$ at $r \neq 0$ and $C_{\alpha\beta}(r) = \langle C_{\alpha\beta, M}(r) \rangle$ we have

$$\langle C_{\alpha\beta, M}^2(r) \rangle = C_{\alpha\beta}^2(r) + \langle C_{\alpha\beta, f}^2(r) \rangle, \quad r \geq 1. \quad (34)$$

The correlation function $C_{\alpha\beta}(r)$ vanishes when r exceeds the correlation length R_c , $r \gg R_c$. It makes possible to find the asymptotical value of $C_{\alpha\beta, f}^2(r)$

$$\langle C_{\alpha\beta, f}^2(r) \rangle|_{r \gg R_c} \cong \langle C_{\alpha\beta, M}^2(r) \rangle =$$

$$\frac{1}{(M-r)^2} \left\langle \sum_{i,j=0}^{M-r-1} [\delta(a_i, \alpha) - p_\alpha] [\delta(a_{i+r}, \beta) - p_\beta] \right.$$

$$\left. \times [\delta(a_j, \alpha) - p_\alpha] [\delta(a_{j+r}, \beta) - p_\beta] \right\rangle. \quad (35)$$

Neglecting the correlations between elements a_i and taking into account that the terms with $i = j$ give the

main contribution to the result,

$$\left\langle \sum_{i,j=0}^{M-r-1} [\delta(a_i, \alpha) - p_\alpha] [\delta(a_{i+r}, \beta) - p_\beta] \right.$$

$$\left. \times [\delta(a_j, \alpha) - p_\alpha] [\delta(a_{j+r}, \beta) - p_\beta] \right\rangle$$

$$\cong \sum_{i=0}^{M-r-1} \langle [\delta(a_i, \alpha) - p_\alpha]^2 \rangle \langle [\delta(a_{i+r}, \beta) - p_\beta]^2 \rangle$$

$$= (M-r) C_{\alpha\alpha, f}(0) C_{\beta\beta, f}(0). \quad (36)$$

we obtain, after neglecting r in the term $M-r$, the averaged fluctuation-dependent contribution to the squared correlation function

$$\langle C_{\alpha\beta, f}^2(r) \rangle \cong \frac{1}{M} C_{\alpha\alpha, f}(0) C_{\beta\beta, f}(0), \quad (37)$$

$$C_{\alpha\beta}(0) = p_\alpha \delta(\alpha, \beta) - p_\alpha p_\beta.$$

Note that Eq. (37) is obtained by means of averaging over the ensemble of chains. This is the shortest way to get the desired result. At the same time, for numerical simulations we have only used the averaging over the chain as is seen from Eq. (32), where the summation over sites i of the chain plays the role of averaging.

Note also that the different symbols a_i in Eq. (36) are correlated. It is possible to show by direct evaluation of $C_{\alpha\beta, f}^2(r)$ with CPF (21) that the contribution of their correlations to $\langle C_{\alpha\beta, f}^2(r) \rangle$ is of order of $\Delta/M^2 \ll 1/M$.

Equation (27), containing $C_{\alpha\beta}(r)$, is only valid for the infinite chain. In reality, we always work with sequences of finite length and can calculate $C_{\alpha\beta, M}(r)$, which contains the fluctuating part. To improve result (27) we have to subtract the fluctuating part of entropy, proportional to $\sum_{r=1}^L \langle C_{\alpha\beta, f}^2(r) \rangle$, from Eq. (27). Thus, Eqs. (34) and (37) yield the differential entropy of the *finite* weakly correlated (approximately ergodic, $R_c \ll M$) random sequences

$$h_L = h_0 - \frac{1}{2 \ln 2} \left[\sum_{r=1}^L \sum_{\alpha, \beta \in A} \frac{C_{\alpha\beta, M}^2(r)}{p_\alpha p_\beta} - (m-1)^2 \frac{L}{M} \right]. \quad (38)$$

It is clear that in the limit $M \rightarrow \infty$ this function transforms into Eq. (27). The last term in RHS of Eq. (38) describes the linearly decreasing fluctuation correction of the entropy. For the binary chain, $m = 2$, we get the result obtained earlier in [10].

The squared correlation function $C_{\alpha\beta, M}^2(r)$ is normally a decreasing function of r , whereas the function $C_{\alpha\beta, f}^2(r)$ is nearly constant (see Eq. (37) for $r \ll M$). Hence, the terms $\sum_{r=1}^L \sum_{\alpha, \beta \in A} C_{\alpha\beta, M}^2(r)/p_\alpha p_\beta$ and $(m-1)^2 L/M$ being concave and linear functions, respectively, describe the competitive contributions to the entropy. It is not possible to analyze all particular cases of their relationship. Therefore we indicate here the most interesting ones taking in mind monotonically decreasing correlation functions. An example of such a type of function is $C(r) = a/r^b$, $a > 0$, $b > 0$. If the correlations are ex-

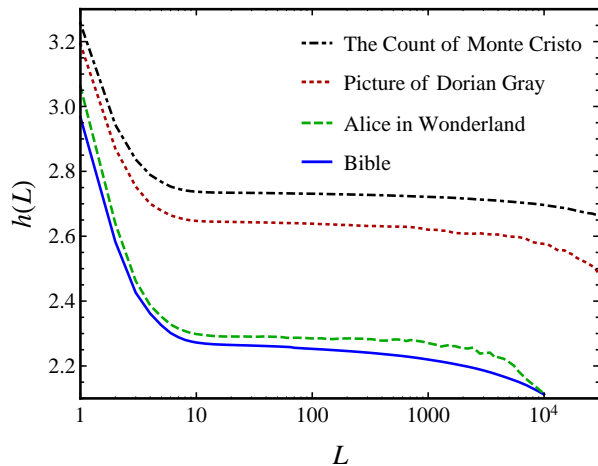


FIG. 1: (Color online) The differential entropy of the literature works (indicated in the legend near the curves) vs the length of words in L -axis log scale. The curves correspond to the direct evaluations of Eq. (27) with fluctuation correction.

tremely small and compared with the inverse length M of the sequence, $\sum_{\alpha, \beta \in A} C_{\alpha\beta, M}^2(1)/p_{\alpha}p_{\beta} \sim 1/M$, the fluctuating part of the entropy exceeds the correlation part almost for all values of $L > 1$.

When the correlations are more strong, $\sum_{\alpha, \beta \in A} C_{\alpha\beta, M}^2(1)/p_{\alpha}p_{\beta} > 1/M$, there is at least one point where the contribution of fluctuation and correlation parts of the entropy are equal. For monotonically decreasing function $\sum_{\alpha, \beta \in A} C_{\alpha\beta, M}^2(r)/p_{\alpha}p_{\beta}$ there is only one such point. Comparing the functions in square brackets in Eq. (38) we find that they are equal at some $L = R_s$, which hereafter will be referred to as a stationarity length. If $L \ll R_s$, the fluctuations of the correlation function are negligibly small with respect to its magnitude, hence for these L -words the finite sequence may be considered as the quasi-stationary one. At $L \sim R_s$ the fluctuations are of the same order as the genuine correlation function contribution, $\sum_{\alpha, \beta \in A} C_{\alpha\beta, M}^2(r)/p_{\alpha}p_{\beta}$. Here we have to take into account the fluctuation correction due to the finiteness of the random chain. At $L > R_s$ the fluctuation contribution exceeds the correlation one and Eq. (38) loses any sense.

The other important parameter of the random sequence is the memory length N . If the length N is less than R_s , we have no difficulties to calculate the entropy of the finite sequence, which can be considered as quasi-stationary. If the memory length exceeds the stationarity length, $R_s \lesssim N$, we should take into account the fluctuation correction to the entropy.

V. APPLICATIONS TO NATURAL AND DNA TEXTS

The purpose of this section is to illustrate applicability of the developed theory to some concrete sequences naturally arising in biology and linguistics.

In order to evaluate the differential entropy of literature works we calculate the probabilities p_{α} of each letter occurring in the simplified text and symbolic correlation functions $C_{\alpha\beta, M}(r)$. The simplification (some sort of coarse-graining) consists in replacing all the upper-case letters with the lower-case ones and neglecting all punctuation marks except blanks. Hence, we use the alphabet of 27 letters. The result for calculating the differential entropy with the use of Eq. (38) is shown in Fig. 1. The entropy per one letter $h(0)$ (not shown in the picture) is 4 ± 0.1 . It is evident that the difference between the one-letter-entropy, in the case of the letters equipartition $\log_2 27 \approx 4.75$, and 4 ± 0.1 is due to the non-equipartition distribution of letters in the texts.

As we mentioned, the correlation length can be determined as the length where the entropy takes on a constant value. At first glance, the value of R_c is of order of 9 – 11. But after this point we observe a nearly linear small decrease of entropy extended over 2 – 3 decades. Probably, this phenomenon could be explained by small power-law correlation observed and discussed in Ref. [17].

Application of the developed theory to nucleotide sequences of DNA molecules is shown in Fig. 2. In order to evaluate the entropy of the *Homo sapiens* chromosome Y, locus NW 001842422 [23], we calculate the probabilities p_{α} of each nucleotide occurring in the sequence and 9 different symbolic correlation functions $C_{\alpha\beta, M}(r)$.

It is clearly seen that the entropy in the interval $7 \times 10^3 < L < 2 \times 10^4$ takes on the constant value, $h_L \simeq 1.41$. It means that for $L > 7 \times 10^3$ all binary correlations, in the statistical sense, are taken into account. In other words, the correlation length of the *Homo sapiens* chromosome Y is of the order of 10^4 . This length R_c is much greater than correlation length $R_c \approx 10$ observed for natural written texts.

In the inset the differential entropy of *Homo sapiens* chromosome Y, locus NW 001842451, is shown. Here we cannot see a constant asymptotical region, which would be an evidence for the existence of stationarity and finiteness of the correlation length. We suppose that the locus is not well described by our theory at long distances due to the relatively short length of sequence. The dashed line in the figure is the fluctuation correction of the differential entropy. This correction should be small with respect to the correlation contribution in the region of reliability of the result. Thus, only for $L < 10^3$ the result can be considered as a plausible.

It is interesting to compare our results with those obtained by estimation of block entropy Eq. (5) where the probabilities of words occurring are calculated with stan-

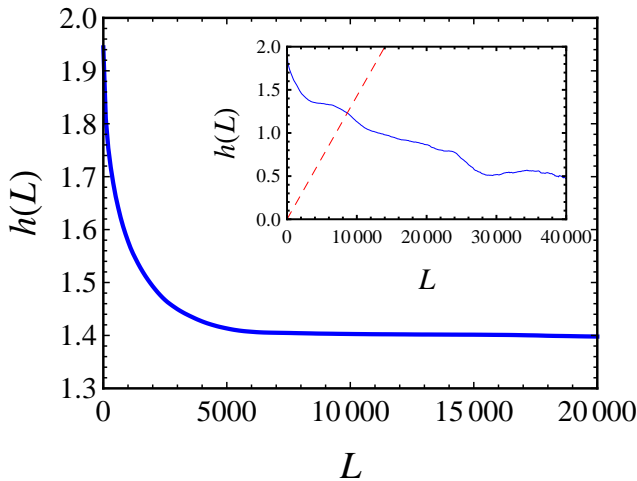


FIG. 2: (Color online) The differential entropy of *Homo sapiens* chromosome Y, locus NW 001842422 [23], of length $M \simeq 3.9 \times 10^6$ vs length L with the fluctuation correction. The curve is constructed by using Eq. (27). The inset demonstrates the differential entropy of *Homo sapiens* chromosome Y, locus NW 001842451, of length $M \simeq 4.5 \times 10^4$. The straight dashed line is fluctuation correction $9L/2 \ln 2M$ due to finiteness of chain.

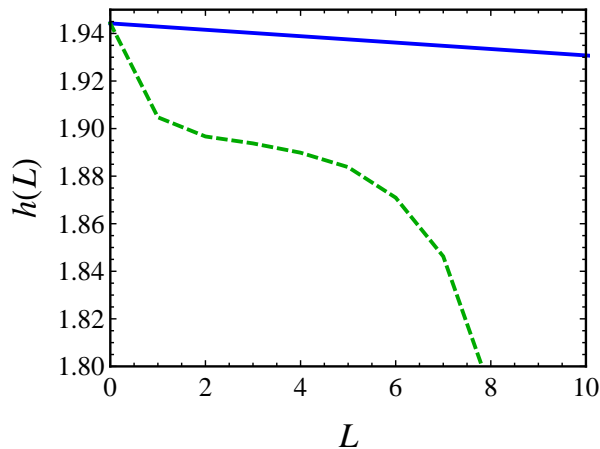


FIG. 3: (Color online) Comparison of differential entropies calculated by estimation of block occurring Eq. (39) (bottom curve) and the result of Eq. (38) (top curve) for the *Homo sapiens* chromosome Y, locus NW 001842422.

dard likelihood estimate

$$P(a_1^L) = \frac{n(a_1^L)}{M - L + 1}. \quad (39)$$

Here $n(a_1^L)$ is the number of occurrences of the word a_1^L in the sequence of the length M . In our paper [10] it was shown that there is a good agreement between two approaches for the coarse-grained (binary) DNA sequence of R3 chromosome of *Drosophila melanogaster* of length $M \simeq 2.7 \times 10^7$ for $L \lesssim 5 - 6$ units. For four-valued sequence (composed by adenine, guanine, cytosine, thymine) we cannot make a similar conclusion

studying the differential entropy of the *Homo sapiens* chromosome Y, locus NW 001842422, shown in Fig. 3. It is clear that at small L strong short-range correlations or the exact statistics of the short words are more important than that which we took into account — the simple pair correlations.

It is difficult to come to an unambiguous conclusion, which factor, the finiteness of the chain and violation of Eq. (1) or the strength of correlations, is more important for the discrepancy between the two theories and between the two studied sequences.

VI. CONCLUSION AND PERSPECTIVES

(i) The main result of the paper, the differential entropy of the stationary ergodic weakly correlated random sequence \mathbb{A} with elements belonging to the finite alphabet is given by Eq. (27). The other important point of the work is the calculation of the fluctuation contribution to the entropy due to the finiteness of random chains, the last term in Eq. (38).

(ii) In order to obtain Eq. (27) we used an assumption that the random sequence of symbols is the high-order Markov chain. Nevertheless, the final result contains only the correlation function and does not contain the conditional probability function of the Markov chain. This allows us to suppose that result (27) and the region of its applicability is wider than the assumptions under which it is obtained.

(iii) To obtain Eq. (27) we supposed that the correlations in the random chain are weak. It is not a very severe restriction. Many examples of such kind of systems described by means of the pair correlator are given in Ref. [1]. The randomly chosen example of DNA sequences and the literary texts support this conclusion. The strongly correlated systems, which is opposed to weakly correlated chains, are nearly deterministic. For their description we need completely different approach. Their study is beyond the scope of this paper.

(iv) Equation (27) can be considered as an expansion of the entropy in series with respect to the small parameter Δ , where the entropy h_0 of the non-correlated sequence is the zero approximation. Alternatively, for the zero approximation we can use the exactly solvable model of the N -step Markov chain with the conditional probability function of words occurring taken in the form of the step-wise function [18]. Another way to choose the zero approximation can be based on CPF obtained from probability of the block occurring Eq. (5). Consequently, the developed theory opens the way to construct a more consistent and sophisticated approach describing the systems with strong short-range and weak long-range correlations.

(v) Our consideration can be generalized to the Markov chain with the infinite memory length N . In this case we should impose the condition of the decreasing rate of the correlation function and the conditional probability

function at $N \rightarrow \infty$.

Appendix A

The conditional probability function of the *binary additive* Markov chain of random variables $a_i \in \{0, 1\}$, the probability of symbol a_i to have a value 1 under the condition that N previous symbols a_{i-N}^{i-1} are given, is of the following form [9, 16],

$$P(a_i = 1 | a_{i-N}^{i-1}) = \bar{a} + \sum_{r=1}^N F(r)(a_{i-r} - \bar{a}). \quad (\text{A1})$$

Analogously for $P(0|\cdot)$,

$$\begin{aligned} P(a_i = 0 | a_{i-N}^{i-1}) &= 1 - P(1 | a_{i-N}^{i-1}) \\ &= 1 - \bar{a} - \sum_{r=1}^N F(r)(a_{i-r} - \bar{a}). \end{aligned} \quad (\text{A2})$$

This two expressions are not symmetric with respect to the change $0 \leftrightarrow 1$ of generated symbol a_i . Let us show that Eqs. (A1) and (A2) can be presented in the symmetric form

$$P(a_i = \alpha | a_{i-N}^{i-1}) = p_\alpha + \sum_{r=1}^N \sum_{\beta \in \{0,1\}} F_{\alpha\beta}(r) [\delta(a_{i-r}, \beta) - p_\beta]. \quad (\text{A3})$$

Taking into account the definitions $p_1 = \bar{a}$, $p_0 = 1 - \bar{a}$, using the evident equalities $\delta(a_{i-r}, 0) = 1 - a_{i-r}$, $\delta(a_{i-r}, 1) = a_{i-r}$ and putting $F_{11}(r) - F_{10}(r) = F_{00}(r) - F_{01}(r) = F(r)$ we easily obtain Eqs. (A1) and (A2). We should replace $\alpha, \beta \in \{0, 1\}$ in Eq. (A3) by $\alpha, \beta \in A$ to obtain Eq. (11).

Note, there is no one-to-one correspondence between the memory function $F_{\alpha\beta}(r)$ and the conditional probability function $P(a_i = \alpha | a_{i-N}^{i-1})$. Indeed, it is easy to see that, in view of Eqs. (11) and (12), the renormalized memory function $F'_{\alpha\beta}(r) = F_{\alpha\beta}(r) + \varphi_\alpha(r)$ provides the same conditional probability as $F_{\alpha\beta}(r)$.

Appendix B

Here we prove Eq. (23) using Eqs. (11) and (21) as a starting point. It follows from definition (10) of the conditional probability function

$$P(a_i = a | W) = \frac{P(W, a)}{P(W)}, \quad W = a_{i-N+1}^{i-1}. \quad (\text{B1})$$

Adding symbol $a_{i-N} = b$ to the string (W, a) we have

$$P(a_i = a | W) = \frac{\sum_{b \in A} P(b, W, a)}{P(W)}. \quad (\text{B2})$$

Replacing here the probabilities $P(b, W, a)$ by the CPF $P(a_i = a | b, W)$ from the equation similar to that of Eq. (B1),

$$P(a_i = a | b, W) = \frac{P(b, W, a)}{P(b, W)}, \quad (\text{B3})$$

we obtain after some algebraic manipulations

$$\begin{aligned} P(a_i = a | W) &= p_a + \sum_{r=1}^{N-1} \sum_{b \in A} F_{ab}(r) [\delta(a_{i-r}, b) - p_b] \\ &+ \frac{1}{P(W)} \sum_{c \in A} F_{ac}(N) \sum_{b \in A} P(b, W) [\delta(b, c) - p_c]. \end{aligned} \quad (\text{B4})$$

The 3-rd term containing summation over b is of the form

$$P(c, W)(1 - p_c) - P(\bar{c}, W)p_c, \quad (\text{B5})$$

where the symbol \bar{c} stands for an event NOT- c . It is intuitively clear that in the zero approximation in Δ (i.e., for uncorrelated sequence) this term equals zero. In the next approximation this term is of order of Δ . These two statements can be verified by using the condition of compatibility for the Chapman-Kolmogorov equation (see, for example, Ref. [24]),

$$P(a_{i-N+1}^i) = \sum_{a_{i-N} \in A} P(a_{i-N}^{i-1}) P_N(a_i | a_{i-N}^{i-1}). \quad (\text{B6})$$

Hence, we have to neglect the third term in the right-hand side of Eq. (B4) because it is of the second order in Δ . So, Eq. (23) is proven for $L = N - 1$. By induction, the equation can be written for arbitrary L .

Acknowledgments

We are grateful for the helpful and fruitful discussions with G. M. Pritula, S. S. Apostolov, and Z. A. Maizelis.

-
- [1] F. M. Izrailev, A. A. Krokhin, N. M. Makarov, Phys. Rep. **512**, 125 (2012).
[2] P. Ehrenfest and T. Ehrenfest, *Encyklopädie der Math-*

- ematischen Wissenschaften* (Springer, Berlin, 1911), p. 742, Bd. II.
[3] D. Lind and B. Marcus. *An Introduction to Symbolic Dy-*

- namics and Coding* (Cambridge University Press, Cambridge, 1995).
- [4] C. E. Shannon and W. Weaver, *The Mathematical Theory of Communication* (University of Illinois Press, Urbana, Illinois, 1949).
 - [5] T. M. Cover, J. A. Thomas, *Elements of Information Theory* (Wiley, New York, 1991).
 - [6] D. Salomon, *A Concise Introduction to Data Compression*, (Springer, Berlin, 2008).
 - [7] C. D. Manning, P. Raghavan and H. Schütze, *Introduction to Information Retrieval*, (Cambridge University Press, Cambridge, 2008).
 - [8] A. D. Wissner-Gross and C. E. Freer, Phys. Rev. Lett. **110**, 168702 (2013).
 - [9] S. S. Melnyk, O. V. Usatenko, V. A. Yampol'skii, Physica A **361**, 405 (2006).
 - [10] S. S. Melnik and O. V. Usatenko, to be published in Phys. Rev. E.
 - [11] A. Raftery, J. R. Stat. Soc. B **47**, 528 (1985).
 - [12] W.K. Ching, E.S. Fung, M.K. Ng, Naval Res. Logist. **51**, 557 (2004).
 - [13] W.K. Li, M.C.O. Kwok, Commun. Stat. Simul. Comput. **19**, 363 (1990).
 - [14] J.A. Cocho, *et al.* Comput. Biol. Chem. **53**, 15 (2014).
 - [15] M. Seifert, A. Gohr, M. Strickert, I. Grosse, PLoS Computat. Biol. **8**, e1002286 (2012).
 - [16] O. V. Usatenko, S. S. Apostolov, Z. A. Mayzelis, and S. S. Melnik, *Random Finite-Valued Dynamical Systems: Additive Markov Chain Approach* (Cambridge Scientific Publisher, Cambridge, 2010).
 - [17] S. S. Melnyk, O. V. Usatenko, V. A. Yampol'skii, V. A. Golick, Phys. Rev. E **72**, 026140 (2005).
 - [18] O. V. Usatenko, V. A. Yampol'skii, Phys. Rev. Lett. **90**, 110601 (2003).
 - [19] A. N. Shiryaev, *Probability* (Springer, New York, 1996).
 - [20] P. Grassberger, arXiv:physics/0207023 [physics.data-an].
 - [21] See, e.g., A. M. Yaglom, *Correlation theory of stationary and related random functions* (Springer-Verlag, New York, 1987).
 - [22] N. Chakravarthy, A. Spanias, L. D. Iasemidis, K. Tsakalis, EURASIP J. Appl. Signal Process. **1**, 13 (2004).
 - [23] <ftp://ftp.ncbi.nih.gov/genomes/>.
 - [24] C. W. Gardiner: *Handbook of Stochastic Methods for Physics, Chemistry, and the Natural Sciences*, Springer Series in Synergetics, Vol. 13 (Springer-Verlag, Berlin, 1985).